

Probabilistic formulation of the classification problem

The goal in a typical classification problem is to take an input x and to assign it to one of K discrete classes C_k where $k = 1, \dots, K$ (given the dataset $D = \{(x_i, c_i)\}_{i=1}^n\}$)
 $f: x \rightarrow C = \{C_1, C_2, \dots, C_K\}$

A probabilistic model does so by computing the posterior probability $P(C_i|x, D)$ = probability that the result is class C_i given a new sample x and the dataset D .

In the most common scenario the classes are taken to be disjoint, so that each input is assigned to one and only one class. The input space is thereby divided into decision regions whose boundaries are called decision boundaries or decision surfaces.

Datasets whose classes can be separated exactly by linear decision surfaces are said to be linearly separable.

We will consider the more general case in which $X = \mathbb{R}^d$ (X is a vector of real numbers): this means that the probability distributions are continuous functions.

There are two families of probabilistic models (two ways to compute $P(C_i|x, D)$) $\leftarrow (P(C_i|x))$ \leftarrow omitting D

- generative models : estimate $P(x|C_i)$ and then compute $P(C_i|x)$ with Bayes
 - discriminative models : estimate $P(C_i|x)$ directly
- \rightarrow class conditional densities

generative models \rightarrow Gaussian Naive Bayes

discriminative models \rightarrow Logistic Regression

Probabilistic Generative Models \rightarrow Gaussian Naive Bayes

- for just two classes:
(concept learning)

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x)} = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

↑ Bayes theorem

$$a = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)} = \frac{1}{1 + e^{-a}} = \sigma(a)$$

\rightarrow definition of the sigmoid function

in the case of two classes we can expand the denominator with the theorem of total probabilities this way



We said that if $X = \mathbb{R}^d$ (continuous vector) and $x \in X$ then the posterior probabilities are continuous functions. We assume that the probability distributions are gaussians.

$$P(x|C_i) = N(x, \mu_i, \Sigma)$$

μ_i mean
 Σ covariance

We also assume that all classes share the same covariance matrix.

Given a vector x , the covariance matrix is a square matrix that gives the covariance between each pair of elements of the vector.

Equivalent formulation for two (or more) classes

	C_1	C_2
C_1	$\sigma_{C_1}^2$	σ_{C_1, C_2}
C_2	σ_{C_2, C_1}	$\sigma_{C_2}^2$

$$a = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)} \longrightarrow a = \ln \frac{N(x, \mu_1, \Sigma)P(C_1)}{N(x, \mu_2, \Sigma)P(C_2)}$$

after a series of steps :

w and w₀ are different for other exponential distributions

$$P(C_1|x) = \sigma(w^T x + w_0)$$

$$w = \Sigma^{-1}(\mu_2 - \mu_1)$$

$$w_0 = -\frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{P(C_1)}{P(C_2)}$$

Under Gaussian assumption $P(C_i|x)$ becomes a linear expression in x (input). This is true even for other distributions that belongs to the exponential family: posterior is a linear combination of the input

this relation holds for exponential distributions

$$a = (-w^T) \begin{pmatrix} 1 \\ x \\ 1 \end{pmatrix} + w_0$$

scalar

$$f: \mathbb{R} \rightarrow \{\mathbb{G}_1, \mathbb{G}_2\}$$

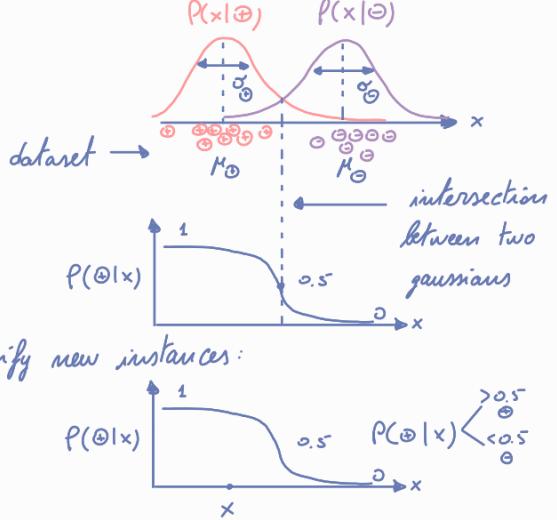
■ multi-class case

$$P(C_k|x) = \frac{P(x|C_k)P(C_k)}{\sum_j P(x|C_j)P(C_j)} = \frac{e^{a_k}}{\sum_j e^{a_j}}$$

called softmax function

$$a_i = \ln P(x|C_k)P(C_k)$$

It is now easy to classify new instances:



Maximum Likelihood

■ for two classes

$$\begin{aligned} P(C_1) &= \pi & P(x|C_1) &= N(x, \mu_1, \Sigma) \\ P(C_2) &= 1-\pi & P(x|C_2) &= N(x, \mu_2, \Sigma) \end{aligned} \quad \left. \begin{array}{l} \text{probabilistic formulation of the concept learning} \\ \text{problem} \end{array} \right\}$$

The parameters are π, μ_1, μ_2 (constraint with $P(C_2) = 1 - \pi$)

Once we have these parameters, we can compute the posterior $P(C_i|x)$ with Bayes

Until now the dataset was always present in the formulas but omitted in order to ease the notation.

$$\Delta = \{(x_m, t_m)\}_{m=1}^N \quad \text{with } t_m = \begin{cases} 1 & \text{if } x_m \in C_1 \\ 0 & \text{if } x_m \in C_2 \end{cases}$$

$$t = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix} \quad P(t | \pi, \mu_1, \mu_2, \Sigma, \Delta) = \prod_{m=1}^N \left[\pi N(x_m, \mu_1, \Sigma) \right]^{t_m} \left[(1-\pi) N(x_m, \mu_2, \Sigma) \right]^{(1-t_m)}$$

if $t_m=1$ the second term is equal to 1 and only the first counts

for each sample, only one of the two gives a contribution

all the samples are independent from each other

Likelihood: probability of the output given the parameters of the model and the dataset

Our problem is to find a maximum likelihood solution: estimate unknown $\pi, \mu_1, \mu_2, \Sigma$ such that the obtained solution is optimal (maximize the likelihood).

$$\underset{\pi, \mu_1, \mu_2, \Sigma}{\operatorname{argmax}} P(t | \pi, \mu_1, \mu_2, \Sigma, \Delta) \longrightarrow \underset{\pi, \mu_1, \mu_2, \Sigma}{\operatorname{argmax}} \log P(t | \pi, \mu_1, \mu_2, \Sigma, \Delta)$$

Solution (no steps): → ML solution is still an approximation (maximum given the dataset)

$$\pi \approx \frac{N_1}{N} \quad \frac{\text{number of examples classified in } C_1}{\text{number of examples}} \quad \frac{|\{x, C_1\}|}{|\Delta|}$$

$$\mu_1 = \frac{1}{N_1} \sum_{m=1}^N t_m x_m \quad \text{with } t_m = \begin{cases} 1 & \text{if } x_m \in C_1 \\ 0 & \text{if } x_m \in C_2 \end{cases} \quad \text{mean of all the samples that belongs to class } C_1$$

this way in μ_1 we only have samples from C_1 and in μ_2 we only have samples from C_2

$$\Sigma = \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2 \quad \text{with } S_i = \frac{1}{N_i} \sum_{m \in C_i} (x_m - \mu_i)(x_m - \mu_i)^T, \quad i=1,2$$

- ## ■ for k classes

$$P(C_1) = \pi_1$$

$$\vdots \quad \vdots$$

$$P(C_N) = 1 - \sum_{i=1}^{N-1} \pi_i$$

$$P(x|C_1) = \mathcal{N}(x, \mu_1, \Sigma)$$

$$P(C_N) = 1 - \sum_{i=1}^{N-1} \pi_i$$

$$p(x|c_u) = \mathcal{N}(x; \mu_u, \Sigma_u)$$

$$\Delta = \left\{ (x_m, t_m) \right\}_{m=1}^N \quad \text{with } t_m \text{ 1-of-k encoding}$$

Solution (no steps) :

$$\pi_i = \frac{N_i}{N}$$

$$\mu_i = \frac{1}{N_i} \sum_{m=1}^N t_{mk} x_m$$

$$\sum = \sum_{k=1}^k \frac{N_k}{N} S_k \quad \text{with} \quad S_k = \frac{1}{N_k} \sum_{m=1}^N \epsilon_{mk} (x_m - M_k) (x_m - M_k)^T$$

one row for each sample of the dataset, one column for each class; In each row we have all zeros except a 1 in the position that corresponds to the class of the sample.

The gaussian assumption does not work well if the classes are not linearly separable

Given a generic model, suppose it has 3 parameters. The solution for the argmax is:

$$\begin{aligned} \frac{\partial \log P(x; \theta_1, \theta_2, \theta_3)}{\partial \theta_1} &= 0 \longrightarrow \theta_1^* \\ \frac{\partial \log P(x; \theta_1, \theta_2, \theta_3)}{\partial \theta_2} &= 0 \longrightarrow \theta_2^* \\ \frac{\partial \log P(x; \theta_1, \theta_2, \theta_3)}{\partial \theta_3} &= 0 \longrightarrow \theta_3^* \end{aligned}$$

Probabilistic Discriminative Models → Logistic Regression

The goal in discriminative models is to directly estimate the posterior without considering the class conditional densities → we can ignore the parameters of the distribution and only use w . Thus, we will not be able to generate new samples.

- for two classes:

$$D = \{(\tilde{x}_m, t_m)_{m=1}^N\} \quad \text{with } t_m \in \{0, 1\}$$

$$P(t|\tilde{w}) = \prod_{m=1}^N y_m^{t_m} (1-y_m)^{1-t_m}$$

likelihood
function

product is an assumption of the independence of the data

prediction of
the posterior

, $y_m = P(C_1 | x_m)$
 prediction that the

prediction that the model with parameters \tilde{w} will give an output one particular instance

(W) will give an input x_m
one particular instance of the

formulate the problem:

$$E(\tilde{w}) = -\ln P(t| \tilde{w}) = -\sum_{m=1}^N [t_m \ln y_m + (1-t_m) \ln (1-y_m)]$$

negative log likelihood

$$\text{problem} = \underset{\omega}{\operatorname{argmin}} E(\tilde{\omega})$$

$$= \underset{\tilde{w}}{\operatorname{argmax}} \ln P(t | \tilde{w})$$

CROSS-ENTROPY (it is an error and is defined as the negative log likelihood)

log likelihood)

We can use any solver able to minimize/maximize a function
 ↳ the Newton-Raphson solver is one of those

The Newton-Raphson method is an iterative optimization method.

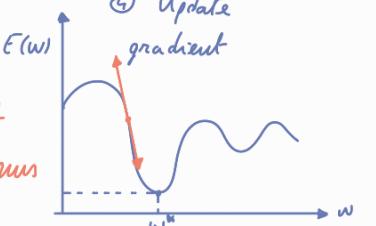
$$\nabla E(\tilde{w}) = \sum_{m=1}^N (y_m - t_m) \tilde{x}_m$$

iteratively:

$$\tilde{w} \leftarrow \tilde{w} - H(\tilde{w})^{-1} \nabla E(\tilde{w})$$

"how much you want to move" computed in terms of the second derivative:

the steeper the function, the greater the update



We move towards the minimum. It is not guaranteed to reach the global minimum; it depends on the starting point.

■ multiclass logistic regression:

$$f: \mathbb{R}^d \rightarrow \{c_1, \dots, c_k\}$$

$$D = \{(x_m, t_m)\}_{m=1}^N \text{ with } t_m = \frac{1}{(0, \dots, 0 \pm 1, 0, \dots, 0)^T} \text{ 1-of-k encoding}$$

$$\tilde{x} = \begin{pmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_N^T \end{pmatrix} \quad T = \begin{pmatrix} t_1^T \\ \vdots \\ t_N^T \end{pmatrix} \quad \text{1-of-k encoding in each row for the labels}$$

\tilde{x} is a $d+1$ dimensional vector, one for each input sample

We have a vector \tilde{w}_i for each class: $\tilde{w}_1, \dots, \tilde{w}_k \quad i \in [1, \dots, k]$

We denote with y_{mk} the prediction of model k on the sample x_m

$y_m^T = (y_{m1}, \dots, y_{mk})^T$ posterior prediction of \tilde{x}_m for model $\tilde{w}_1, \dots, \tilde{w}_k$

$y_{mk} = \frac{\tilde{w}_k^T \tilde{x}_m}{\sum_j \tilde{w}_j^T \tilde{x}_m}$ one prediction for each class

$$a_k = \tilde{w}_k^T \tilde{x} \quad k \in [1, \dots, k] \rightarrow P(c_k | \tilde{x}) = \frac{e^{a_k}}{\sum_i e^{a_i}}$$

$$Y(\tilde{w}_1, \dots, \tilde{w}_k) = \begin{pmatrix} y_1^T \\ \vdots \\ y_N^T \end{pmatrix} \quad \text{each row is a prediction vector for the respective sample: a vector of K components, one for each class. Sum of each row is = 1}$$

Once we defined this, we can define the likelihood for the K classes case:

$$P(T | \tilde{w}_1, \dots, \tilde{w}_k) = \prod_{m=1}^N \prod_{k=1}^K P(c_k | \tilde{x}_m)^{t_{mk}} = \prod_{m=1}^N \prod_{k=1}^K y_{mk}^{t_{mk}}$$

with $y_{mk} = Y[m][k]$ and $t_{mk} = T[m][k]$

Cross entropy for K classes:

$$E(\tilde{w}_1, \dots, \tilde{w}_k) = -\ln P(T | \tilde{w}_1, \dots, \tilde{w}_k) = -\sum_{m=1}^N \sum_{k=1}^K t_{mk} \ln y_{mk}$$

Iterative algorithm:

$$\text{gradient: } \nabla_{\tilde{w}_k} E(\tilde{w}_1, \dots, \tilde{w}_k)$$

$$\text{hessian: } \nabla_{\tilde{w}_k} \nabla_{\tilde{w}_j} E(\tilde{w}_1, \dots, \tilde{w}_k)$$

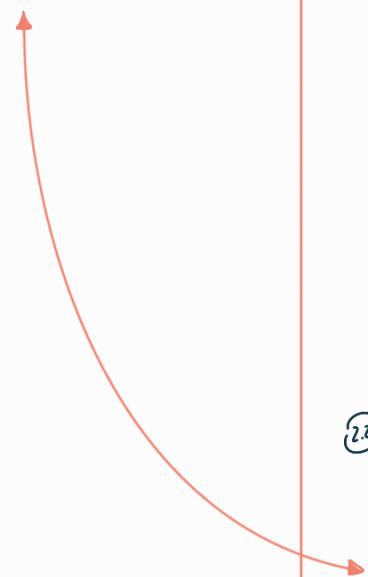
$$\text{iterative step: } \tilde{w}_k \leftarrow \tilde{w}_k - H(\tilde{w}_1, \dots, \tilde{w}_k)^{-1} \nabla E(\tilde{w}_1, \dots, \tilde{w}_k)$$

- ① Random value
- ② Compute gradient
- ③ Invert and follow it
- ④ Update

Generalization

- ① We are given a target function $f: X \rightarrow C$ and a dataset D
- ② Assume a parametric prediction model $y(x, \theta)$ where θ are the parameters of the model and x is the input. y needs to be as close as possible to the target function f

$$y(x, \theta) \approx f(x)$$



Probabilistic models

- ① We are given a target function $f: X \rightarrow C$ and a dataset D
- ② Assume a parametric model for the posterior probability $P(C_k | \tilde{X}, \tilde{\omega})$
- ③ Define a posterior probability $P(C_k | \tilde{x})$
 - $y_m = P(C_1 | \tilde{x}_m)$ for 2 classes
 - $y_{mk} = \frac{e^{\tilde{\omega}_k^T \tilde{x}_m}}{\sum_j e^{\tilde{\omega}_j^T \tilde{x}_m}}$ for K classes

Posterior probability = probability of having class C given the sample and the dataset (after collecting the data)
- ④ Define the likelihood function
 - $P(t | \tilde{\omega}) = \prod_{m=1}^N y_m^{t_m} (1-y_m)^{1-t_m}$ for 2 classes
 - $P(T | \tilde{\omega}_1, \dots, \tilde{\omega}_K) = \prod_{m=1}^N \prod_{k=1}^K P(C_k | \tilde{x}_m)^{t_{mk}}$
 $= \prod_{m=1}^N \prod_{k=1}^K y_{mk}^{t_{mk}}$ for K classes

- ③ Define an error function $E(\theta)$

- ③ Define an error function $E(\tilde{\omega})$

$$E(\tilde{\omega}) = -\ln P(C_k | \tilde{x}, \tilde{\omega})$$

- ④ Solve the optimization problem

$$\theta = \underset{\theta}{\operatorname{argmin}} E(\theta)$$

- ④ Solve the optimization problem

$$\tilde{\omega}^* = \underset{\tilde{\omega}}{\operatorname{argmin}} E(\tilde{\omega})$$

The methods described above (Gaussian Naive Bayes + generalization) can be applied in a transformed space of the input (feature space). Given a function $\phi: \tilde{X} \rightarrow \phi$ (feature space), each sample \tilde{x}_m can be mapped to a feature vector $\phi_m = \phi(\tilde{x}_m)$.

NB: regression = solution of a problem in which the output function is continuous;
logistic regression = classification method.